

# On Seeking Consensus Between Document Similarity Measures

Mieczysław Kłopotek

Institute of Computer Science, Polish Academy of Sciences  
ul. Jana Kazimierza 5, 01-248 Warszawa, Poland  
kłopotek@ipipan.waw.pl

February 14, 2017

## Abstract

This paper investigates the application of consensus clustering and meta-clustering to the set of all possible partitions of a data set. We show that when using a "complement" of Rand Index as a measure of cluster similarity, the total-separation partition, putting each element in a separate set, is chosen.

**Index Terms**— cluster analysis, partitioning, clustering, consensus functions, ensemble, knowledge reuse, unsupervised learning, meta-clustering

## 1 Introduction

It is a well-known phenomenon that for the same data set various clustering algorithms may produce different partitions. This is true both for objects described by continuous variables (like results of measurements) and for ones described by discrete features (like documents treated as points in term space). Consensus clustering and meta-clustering are two known techniques helping to select the best one among the competing partitions. It is also well known that by changing the geometry of the data space we may even obtain all possible partitions of the dataset.

In this paper we investigate which partition would be selected if we apply consensus clustering or meta-clustering to the set of all possible partitions. In particular we formulate (in Section 3) and prove (in Section 4) that, using the so-called Rand Index as a measure of partition similarity, we obtain via consensus clustering the partition putting each element in a separate set (which we will call subsequently *total-separation partition*) as a "consensus" between these partitions. In Section 5 we discuss briefly practical lessons from this theorem. In Section 6 we demonstrate that a similar theorem can be formulated for the more realistic case where we consider only partitions containing not more

clusters than a predefined threshold. In Section 7 we show experimentally that also the very same similarity measure applied in meta-clustering<sup>1</sup> leads towards a similar choice of best partition.

We start with Section 2 explaining the concepts of consensus clustering and meta-clustering as well as pointing to the research on these topics. In Section 8 we summarise our findings and point to further research directions.

## 2 Previous work

Two main types of methods for handling the grieving issue of conflicting partitions of the same set are currently under development in the literature:

- the meta-clustering
- the consensus clustering (called also in various brands ensemble clustering [27] or cluster aggregation [10])

In the *meta-clustering* stream it is claimed that maybe the choice of the best partition should be left to the user who should only be assisted by grouping potential partitions into groups of similar ones. To facilitate user selection of the right clustering, [5] (also compare [21, 4, 3, 7, 6]) suggests to provide the user with meta-clusters (clusters of partitions) in order that the user better understands the choices. To facilitate their creation, [5] proposes to use a dissimilarity measure that they call "Cluster Difference", closely related to Rand Index (which is a similarity measure) as distance measure between partitions. The sum of Rand Index and Cluster Difference is equal 1. The Rand Index (and as a consequence Cluster Difference) essentially is based on the calculation how many times elements are in the same or in different clusters. Assume that the set  $\mathbb{X}$  to be clustered, of cardinality  $n$  consists of elements  $\{1, 2, \dots, n\}$ . Let the quantity  $I_{ij}$  be equal to 1 if the elements (objects)  $i$  and  $j$  are in the same cluster of the first clustering, and different in the second one or vice versa. Otherwise,  $I_{ij}$  is 0. Then the distance  $CD$  between the two partitions  $\Gamma_1, \Gamma_2$  (Cluster Difference) is defined as

$$CD(\Gamma_1, \Gamma_2) = \frac{\sum_{i,j \in \mathbb{X}, i < j} I_{ij}}{n(n-1)/2} \quad (1)$$

where  $n$  is the number of objects in the collection. The value of  $CD$  ranges from 0 for completely identical partitions to 1 for extremely different ones. The extremes are e.g. two partitions: all-in-one ( $\Gamma$  consisting of exactly one cluster containing all elements) and total-separation (every element in a separate cluster).

---

<sup>1</sup>Rand Index is frequently used in both consensus clustering and meta-clustering in search for a compromise clustering

If we wish to compare only partitions over the same set of elements (with cardinality  $n$ ), we can use the unnormalised version of  $CD$ :

$$unCD(\Gamma_1, \Gamma_2) = \sum_{i,j \in \mathbb{X}, i < j} I_{ij} \quad (2)$$

In such a case the  $unCD$  between all-in-one and total-separation partitions amounts to  $\frac{n(n-1)}{2}$ .

But here we encounter the problem:  $n$  objects can be divided into  $k$  clusters in  $O(k^n)$  ways<sup>2</sup>. So we have to do with an NP-hard task.

In *consensus clustering* [24] a kind of optimisation problem (combinatorial optimisation) is formulated and solved. A similarity measure between partitions is introduced and data is re-clustered to get a clustering close to the original ones. Alternatively groups of clusters are formed (a kind of meta-clustering) where the meta-clusters compete for objects performing just a re-clustering. A number of other techniques in this direction was reviewed in [11, 14, 9, 16, 22, 18, 25, 20, 28, 26]. We are particularly interested in the one initiated by [2, 13], and further studied [12] and applied [19]. This approach to consensus clustering seeks to find a new clustering that is as close as possible to all the partitions obtained. The closeness is estimated e.g. (again) as the averaged Rand Index [19].

It has been noticed in the past that various clustering quality indices are biased, and in particular the Rand Index which we discuss in this paper, see for example [15, 17, 8]. It has been reported that Rand Index tends to prefer smaller clusters rather than bigger ones when using it as an external cluster validity index, as for example in the study [8], which concentrated on clusterings into the same number of clusters. [8] demonstrates both theoretically and empirically, that with the increase of the number of clusters, Rand index quickly heads towards stating that compared clusterings are identical. However, in this paper we do not handle the case of clusterings into one fixed number of clusters but rather allow for any number of clusters in a partition.

[15] demonstrates, when studying balanced partitions with different number of clusters, that this time Rand index behaves differently, thus invalidating the generalisations from [8]. The direction of the bias depends on the ground truth clustering (see Theorem 1 in [15] and further ones). The reversal of the bias was also reported from empirical studies [17]. So in fact, the bias tendency of Rand Index remains under these circumstances undecided if we do not know the ground truth partition or we do not know whether it is balanced or not.<sup>3</sup> Note that the paper [15] and other concentrate on the relationship between an empirical partition and the ground truth.

---

<sup>2</sup> More precisely, [1] shows that this number amounts to

$$\frac{1}{k!} \sum_{j=1}^k (-1)^{k-j} \binom{k}{j} j^n \quad (3)$$

<sup>3</sup> Note that the bias has also been studied in the context of stable level of the indicator for randomly assigned partitions under various conditions, see e.g. [23].

However, a study on Rand index bias in the context of consensus clustering and meta-clustering seems not have been performed so far. In such a case, the ground truth cannot be referred to because it is not available.<sup>4</sup> Furthermore, the previous studies concentrate on "preferences" of Rand Index without questioning the existence of ground truth. We demonstrate here that it points at a clustering even if there is no ground truth, no really discernible base clustering.

### 3 Theorem on consensus clustering

It is well known that when we treat all possible geometries of the data set, we can obtain all possible partitions of the data set. Which one shall we choose? Under these circumstances, as we will show, consensus clustering is not useful in the selection process of the partition that is the closest to all the other ones. Because the closest partition is a partition that puts each element (object, document) into a separate cluster.

In particular we will prove the following

**Theorem 1.** *For any  $n$  objects, among all partitions the partition where each object falls into a separate cluster (called subsequently total-separation partition) has the lowest average distance to the other partitions in terms of Cluster Difference CD.*

As we consider a fixed  $n$ , let us concentrate on the unnormalised version  $unCD$ .

Note first one of the most serious implications of this theorem: We consider a world of all possible partitions so one might think that this world is totally symmetric, and any partition may be a centre of such a universe. But this is not the case. The distance function distinguishes one of the elements. So in fact it is biased in some way.

There exist plenty of other clustering quality assessment functions and a similar analysis should be performed for them.

So let us state beforehand that when performing a clustering task, we shall pick at random neither the clustering function, nor the distance function nor the quality assessment function because each of them is biased and we shall care whether or not each function reflects our business purposes.

### 4 The proof

It is assumed that all the elements (objects, documents) of a set to be partitioned possess identifiers being consecutive numbers starting with 1. The proof will be performed by induction on relabelling the objects in a cyclic manner combined with narrowing the set of candidates for the closest element.

First step of checking validity of the theorem (subsection 4.1) for small  $n = 2, 3$  is trivial, but still necessary. Subsequent subsections seek to establish the

---

<sup>4</sup>One may say that we use Rand Index rather as internal and not external quality measure.

induction step by demonstrating a special role of the so-called simple extension of a partition (to be introduced in subsection 4.2, along with the concept of reduct). The important feature here is that total-separation partition of  $n + 1$  elements is a simple extension of total-separation partition of  $n$  elements. Furthermore, as demonstrated in subsection 4.3, distances between partitions of  $n + 1$  elements can be derived from distances between their reducts.

The idea of the inductive step is as follows. First we consider all extensions of a single partition. It is shown in subsection 4.4 that on average the simple extension is the closest one to all the other extensions.

Then in subsection 4.5 we establish, that among all extensions of one partition the simple extension is on average the closest one to all the extensions of some other partition.

These two facts mean that among all extensions of a partition of a set of  $n$  elements the simple extension is on average the closest one to all the partitions of a set of  $n + 1$  elements. Hence, when looking for the consensus partition among all partitions of a set of  $n + 1$  elements we need to consider only those partitions that are simple extensions of partitions of  $n$  elements. In subsection 4.6 we prove that among these candidates (the simple extensions) the simple extension of total-separation partition of  $n$  elements, that is total-separation partition of  $n + 1$  elements is the closest one on average. The induction proceed as follows. First we establish that the solution (the partition closest to all) is among those partitions that have the  $(n + 1)$ st element in a singleton cluster and at least 0 first elements being in singleton clusters. Then we have the inductive step: If the solution is among those partitions that have the  $(n + 1)$ st element in a singleton cluster and at least  $i$  first elements being in singleton clusters, then the solution is among those partitions that have the  $(n + 1)$ st element in a singleton cluster and at least  $i + 1$  first elements being in singleton clusters. After  $n$  inductive steps all  $n + 1$  elements will be in singleton clusters that is we get the total-separation partition as the solution.

#### 4.1 Cases $n = 2, 3$

Consider the unnormalised version of Cluster Difference.

If  $n = 2$ , then there exist only two partitions:  $\Gamma_{1;2} = \{\{1\}, \{2\}\}$  and  $\Gamma_{2;2} = \{\{1, 2\}\}$ . The *unCD* (as well as *CD*) between them equals 1. So for both the average is identical and minimal. The theorem is O.K.

With  $n=3$  we get partitions

- $\Gamma_{1;3} = \{\{1\}, \{2\}, \{3\}\}$  (average *unCD* distance to other partitions **1.5**, (normalised *CD* 0.5),
- $\Gamma_{2;3} = \{\{1, 3\}, \{2\}\}$  (average *unCD* distance to other partitions 1.75)
- $\Gamma_{3;3} = \{\{1\}, \{2, 3\}\}$  (average *unCD* distance to other partitions 1.75)
- $\Gamma_{4;3} = \{\{1, 2\}, \{3\}\}$  (average *unCD* distance to other partitions 1.75)
- $\Gamma_{5;3} = \{\{1, 2, 3\}\}$  (average *unCD* distance to other partitions 2.25).

Theorem is also in this case O.K.

## 4.2 Case $n \rightarrow n + 1$ - reducts and extensions

Consider now what happens when we have computed the unnormalised Cluster Difference between partitions for  $n$  elements and want to compute it for  $n + 1$  elements.

Each partition  $\Gamma$  of  $n + 1$  elements has a unique partition  $\Gamma^*$  with  $n$  elements (called its reduct) from which it can be derived by adding the  $(n + 1)$ st element to an existent cluster or by forming a new one.  $\Gamma$  on the other hand is called an *extension* of  $\Gamma^*$

Let us introduce the concepts of extension and reduct more formally.

**Definition 1.** *Let  $\Gamma^*$  be a partition of  $n$  elements.*

- *If  $\Gamma$  is a partition of  $n + 1$  elements such that  $\Gamma = \Gamma^* \cup \{\{n + 1\}\}$ , then  $\Gamma$  will be called a simple extension of  $\Gamma^*$  ( $\Gamma = \text{simpleextension}(\Gamma^*)$ )*
- *If  $\Gamma$  is a partition of  $n + 1$  elements such that there exists a set  $S \in \Gamma^*$  such that  $\Gamma = (\Gamma^* - \{S\}) \cup \{S \cup \{n + 1\}\}$ , then  $\Gamma$  will be called a complex extension of  $\Gamma^*$*

*Both simple extension and complex extension are extensions. In both cases  $\Gamma^*$  will be called a reduct of  $\Gamma$  ( $\Gamma^* = \text{reduct}(\Gamma)$ ). With  $\text{allextensions}(\Gamma^*)$  we shall denote the set of all (simple and complex) extensions of  $\Gamma^*$ .*

We distinguish complex and simple extensions in order to emphasise the role of a simple extension among the extensions of a partition – on the one hand due to the simplicity of derivation of distances between extensions from the distances between their reducts (subsection 4.3), and on the other because a simple extension is closer to all the other extensions of the same partition than any complex extension, as will be shown later in Section 4.4.

**Example 1.** *It is easily seen, using the notation from the previous section that  $\Gamma_{1;3}, \Gamma_{2;3}, \Gamma_{3;3}$  are extensions of  $\Gamma_{1;2}$ , while  $\Gamma_{4;3}, \Gamma_{5;3}$  are extensions of  $\Gamma_{2;2}$ .  $\Gamma_{1;3}, \Gamma_{4;3}$  are hereby simple extensions, while  $\Gamma_{2;3}, \Gamma_{3;3}, \Gamma_{5;3}$  are complex extensions.*

For another example of a reduct, simple and complex extensions see Figure 1.

## 4.3 Distances between partitions and their reducts

So consider partitions  $\Gamma_1, \Gamma_2$  of  $n + 1$  elements, being extensions of two partitions  $\Gamma_1^*, \Gamma_2^*$  of  $n$  elements resp. Let us denote with  $\text{unCD}(\Gamma_1^*, \Gamma_2^*; n)$  the unnormalised Cluster Difference  $\text{unCD}(\Gamma_1^*, \Gamma_2^*)$ , where the parameter  $n$  draws our attention to the fact that both partitions are defined over the set  $\{1, \dots, n\}$ .

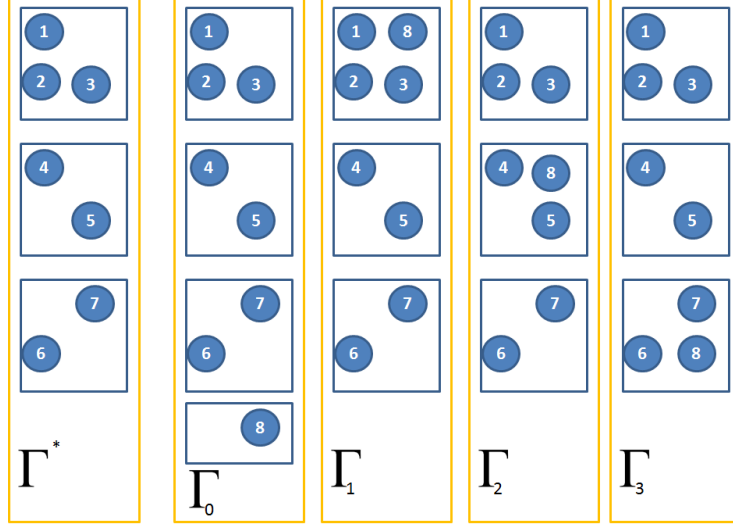


Figure 1: Illustration of the concept of reducts, simple and complex extensions. The partition  $\Gamma^* = \{\{1, 2, 3\}\{4, 5\}\{6, 7\}\}$  is a reduct for each of the partitions  $\Gamma_0 = \{\{1, 2, 3\}\{4, 5\}\{6, 7\}\{8\}\}$ ,  $\Gamma_1 = \{\{1, 2, 3, 8\}\{4, 5\}\{6, 7\}\}$ ,  $\Gamma_2 = \{\{1, 2, 3\}\{4, 5, 8\}\{6, 7\}\}$  and  $\Gamma_3 = \{\{1, 2, 3\}\{4, 5\}\{6, 7, 8\}\}$ .  $\Gamma_1, \Gamma_2, \Gamma_3$  are complex extensions of  $\Gamma^*$ .  $\Gamma_0$  is a simple extension of  $\Gamma^*$ .

Note that

$$\begin{aligned}
 unCD(\Gamma_1, \Gamma_2; n+1) &= \sum_{i=1}^n \sum_{j=i+1}^{n+1} I_{ij} \\
 &= \sum_{i=1}^n \sum_{j=i+1}^n I_{ij} + \sum_{i=1}^n I_{i, n+1} \\
 &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n I_{ij} + \sum_{i=1}^n I_{i, n+1} \\
 &= unCD(\Gamma_1^*, \Gamma_2^*; n) + \sum_{i=1}^n I_{i, n+1}
 \end{aligned} \tag{4}$$

This implies that if both  $\Gamma_1, \Gamma_2$  are simple extensions of  $\Gamma_1^*, \Gamma_2^*$  that is  $\Gamma_1 = \Gamma_1^* \cup \{\{n+1\}\}$  and  $\Gamma_2 = \Gamma_2^* \cup \{\{n+1\}\}$ , then the unnormalised Cluster Difference between  $\Gamma_1$  and  $\Gamma_2$  is the same as between  $\Gamma_1^*$  and  $\Gamma_2^*$  because both in  $\Gamma_1$  and  $\Gamma_2$  element  $n+1$  is separated from any other element.

If  $\Gamma_1 = \Gamma_1^* \cup \{\{n+1\}\}$  and  $\Gamma_2^* - \Gamma_2 = \{S_2\}$  where  $S_2$  is not empty, then  $unCD(\Gamma_1, \Gamma_2; n+1) = unCD(\Gamma_1^*, \Gamma_2^*; n) + card(S_2)$ .

If  $\Gamma_1^* - \Gamma_1 = \{S_1\}$  and  $\Gamma_2^* - \Gamma_2 = \{S_2\}$  where  $S_1, S_2$  are both not empty, then the distance  $unCD(\Gamma_1, \Gamma_2; n+1) = unCD(\Gamma_1^*, \Gamma_2^*; n) + card(S_2 - S_1) +$

$\text{card}(S_1 - S_2)$ .

#### 4.4 Centricity of a simple extension among all extensions

In this proof extensions of a partition play a very special role, because they constitute units for which properties of cumulative distances can be derived in closed form. In particular we show in this subsection that among the extensions of a partition the simple extension is the closest one to all the other extensions. In the next subsection we will demonstrate a similar property between the extensions of two different partitions. The formulas of this subsection are in fact special cases of those in the next subsection, but we believe that by separation of these cases the derivations will be easier to understand.

Let us now consider all extensions with  $n + 1$  elements of a partition  $\Gamma^*$  of  $n$  elements. Let  $\Gamma^*$  contain  $k$  clusters  $S_1, \dots, S_k$ .  $\Gamma_0$  be the simple extension and  $\Gamma_l$  be a complex extension containing the cluster  $S_l \cup \{n + 1\}$ .

Let us compute the sum of distances of simple extension to all the other extensions.

$$\sum_{l=1}^k \text{unCD}(\Gamma_0, \Gamma_l) = \sum_{l=1}^k \text{card}(S_l) = n \quad (5)$$

Let us derive the formula for the sum of distances of a complex extension  $\Gamma_{l'}$  to all the other extensions (remember that clusters of a partition are disjoint,  $S_0$  be an empty set).

$$\begin{aligned} \sum_{l=0, l \neq l'}^k \text{unCD}(\Gamma_l, \Gamma_{l'}) &= \sum_{l=0, l \neq l'}^k (\text{card}(S_l - S_{l'}) + \text{card}(S_{l'} - S_l)) \\ &= \sum_{l=0, l \neq l'}^k (\text{card}(S_l) + \text{card}(S_{l'})) \\ &= k \cdot \text{card}(S_{l'}) + (n - \text{card}(S_{l'})) \\ &= (k - 1) \cdot \text{card}(S_{l'}) + n \geq n \end{aligned} \quad (6)$$

Obviously, as the cardinality of the set of extensions is fixed, the simple extension has the lowest average distance to other extensions among extensions of the same reduct.

**Example 2.** Consider the partition  $\Gamma^* = \{\{1, 3\}\{2\}\}$  and its extensions. The sums of distances of each of them to all the other are: 5 for  $\{\{1, 3, 4\}\{2\}\}$ , 4 for  $\{\{1, 3\}\{2, 4\}\}$ , 3 for  $\{\{1, 3\}\{2\}\{4\}\}$ . The last one is the simple extension and has the lowest sum of distances.

#### 4.5 Distance from a simple and a complex extension

Let us now consider a simple extension  $\Gamma_0$  and a complex one  $\Gamma_m$ ,  $m > 0$  having the same reduct  $\Gamma^*$  and all the partitions  $\Gamma'_l$  with a different common reduct  $\Gamma^{*'}$ . Assume that both  $\Gamma^*$  and  $\Gamma^{*'}$  are partitions over the set  $\{1, \dots, n\}$ .



Let  $\Gamma^{*'} contain  $k$  clusters  $S'_1, \dots, S'_k$ .  $\Gamma_0$  be the simple extension and  $\Gamma_l$  be a complex extension containing the cluster  $S_l \cup \{n+1\}$ .  $S_l$  defined as previously$

Let us calculate the sum of distances of simple extension  $\Gamma_0$  to all the extensions of  $\Gamma^{*'}$ .

$$\begin{aligned} \sum_{l=0}^k unCD(\Gamma_0, \Gamma'_l) &= \sum_{l=0}^k (unCD(\Gamma^*, \Gamma^{*'}; n) + card(S'_l)) \\ &= n + \sum_{l=0}^k unCD(\Gamma^*, \Gamma^{*'}; n) \\ &= n + (k+1)unCD(\Gamma^*, \Gamma^{*'}; n) \end{aligned} \quad (7)$$

Let us determine the sum of distances of complex extension  $\Gamma_m$  to all the extensions of  $\Gamma^{*'}$ .

$$\begin{aligned} \sum_{l=0}^k unCD(\Gamma_m, \Gamma'_l) &= \sum_{l=0}^k (unCD(\Gamma^*, \Gamma^{*'}; n) + card(S_m - S'_l) + card(S'_l - S_m)) \\ &= \sum_{l=0}^k (unCD(\Gamma^*, \Gamma^{*'}; n) + card(S_m) - card(S_m \cap S'_l) \\ &\quad + card(S'_l) - card(S'_l \cap S_m)) \\ &= \sum_{l=0}^k unCD(\Gamma^*, \Gamma^{*'}; n) + (k+1) \cdot card(S_m) + n \\ &\quad - \sum_{l=0}^k card(S_m \cap S'_l) - \sum_{l=0}^k card(S'_l \cap S_m) \\ &= (k+1)unCD(\Gamma^*, \Gamma^{*'}; n) + (k+1) \cdot card(S_m) + n \\ &\quad - card(S_m) - card(S_m) \\ &= (k+1)unCD(\Gamma^*, \Gamma^{*'}; n) + (k-1) \cdot card(S_m) + n \\ &\geq (k+1)unCD(\Gamma^*, \Gamma^{*'}; n) + n \end{aligned} \quad (8)$$

Obviously, as the number of elements in the set of extensions is fixed, the simple extension of  $\Gamma^*$  has the lowest average distance to those extensions of  $\Gamma^{*'}$  among extensions of the  $\Gamma^*$ .

**Example 3.** Consider the partitions  $\mathbf{\Gamma} = allextensions(\{\{1, 3\}\{2\}\})$  and let  $\mathbf{\Gamma}' = allextensions(\{\{1, 2\}\{3\}\})$ . Let us compute for each partition from the set  $\mathbf{\Gamma}$  the sum of distances to all partitions from  $\mathbf{\Gamma}'$ . We obtain: 11 for  $\{\{1, 3, 4\}\{2\}\}$ , 10 for  $\{\{1, 3\}\{2, 4\}\}$ , 9 for  $\{\{1, 3\}\{2\}\{4\}\}$ . The last partition is the simple extension and has the lowest sum of distances.

Therefore we can conclude that if among the extensions of  $\Gamma^*$  there exists a partition that is on average the closest one to any other partition, then this partition is for sure the simple extension of  $\Gamma^*$ .

## 4.6 The partition closest to all the others

So we can summarize sections 4.4 and 4.5 as follows:

**Lemma 1.** *A partition of  $n + 1$  elements closest (on average) to all the other partitions is among the simple extensions of all the partitions of  $n$  elements.*

We can strengthen Lemma 1 by stating:

**Lemma 2.** *If among extensions of partitions of  $n$  elements in the set  $\Gamma$  there is a partition of  $n + 1$  elements closest (on average) to all the other partitions of  $n + 1$  elements, then such a partition exists among simple extensions of  $\Gamma$ .*

Now we are ready to prove by induction that the total-separation partition is the closest on average to all partitions.

Our working hypothesis is as follow: For any  $i = 0, \dots, n$ , the solution of the problem of finding the partition of  $n + 1$  elements closest on average to all the possible partitions of the same set is contained in the set  $C_i$  of such partitions for which the  $(n + 1)$ st element constitutes a singleton in this partition (a cluster with one element only) and the elements  $1, 2, \dots, i$  constitute also singletons. Obviously, if  $i = n - 1$ , then the set containing the solution contains one element only that is the total-separation partition (if all elements but  $n$  are singletons, then also  $n$  is).

Let us first establish the validity of the claim for  $i = 0$ . Let  $\Gamma_{A(n)}$  denote the set of all partitions of  $n$  elements. Lemma 1 states that a partition closest on average to all partitions in  $\Gamma_{A(n+1)}$  is among simple extensions of  $\Gamma_{A(n)}$ . So this is exactly the set of candidates for the on average closest elements denoted previously as  $C_0$ . So initially  $C_0 = \text{simpleextension}(\Gamma_{A(n)})$ , where  $\text{simpleextension}()$  is a function producing simple extensions of partitions.

Next we shall prove the inductive step. That is that if  $C_i$  contains the solution then  $C_{i+1}$  contains also the solution. For this purpose consider the operation of re-labelling elements of partitions. If  $\Gamma$  is a partition, then  $\text{relabel}(\Gamma)$  is a partition obtained by changing an identifier  $i$  of an element to  $i + 1$  except for the element with the highest identifier  $n + 1$  that will be turned to 1. It is obvious that  $\text{unCD}(\Gamma_1, \Gamma_2) = \text{unCD}(\text{relabel}(\Gamma_1), \text{relabel}(\Gamma_2))$ . It is also obvious that  $\text{relabel}(\Gamma_{A(n+1)}) = \Gamma_{A(n+1)}$ , though in general  $\text{relabel}(C_i) \neq C_i$ . However, an element closest on average to each element of  $\Gamma_{A(n+1)}$  is among  $\text{relabel}(C_i)$ . But note that  $\text{simpleextension}(\text{reduct}(\text{relabel}(C_i))) \subseteq \text{relabel}(C_i)$  (as we consider  $n \geq 3$ ). Therefore according to Lemma 2 we can obtain a new candidate set by the operation  $C_{i+1} := \text{simpleextension}(\text{reduct}(\text{relabel}(C_i)))$ . This is because in  $\text{relabel}(C_i)$  in each partition elements  $1, \dots, i + 1$  are singletons, as  $1, \dots, i$  were singletons in  $C_i$  as well as  $n + 1$  and now via relabeling they became  $2, \dots, i + 1$  and 1 respectively. The operation  $\text{simpleextension}(\text{reduct}())$  eliminates everything from  $\text{relabel}(C_i)$  except for simple extensions which have  $n + 1$  as singletons. This proves the validity of the induction step. Obviously, the set of candidates will be reduced in this way.

By induction our claim is valid. Theorem 1 is proven.

**Example 4.** Consider the set

$$\begin{aligned}\Gamma_{A(4)} = & \{ \{ \{1, 2, 3, 4\} \} , \{ \{1, 2, 3\} \{4\} \} , \{ \{1, 2, 4\} \{3\} \} , \{ \{1, 2\} \{3, 4\} \} , \\ & \{ \{1, 2\} \{3\} \{4\} \} , \{ \{1, 3, 4\} \{2\} \} , \{ \{1, 3\} \{2, 4\} \} , \{ \{1, 3\} \{2\} \{4\} \} , \\ & \{ \{1, 4\} \{2, 3\} \} , \{ \{1\} \{2, 3, 4\} \} , \{ \{1\} \{2, 3\} \{4\} \} , \{ \{1, 4\} \{2\} \{3\} \} , \\ & \{ \{1\} \{2, 4\} \{3\} \} , \{ \{1\} \{2\} \{3, 4\} \} , \{ \{1\} \{2\} \{3\} \{4\} \} \}\end{aligned}$$

The set of candidates  $C_0$ , being the simple extensions among the above, is

$$\begin{aligned}C_0 = & \{ \{ \{1, 2, 3\} \{4\} \} , \{ \{1, 2\} \{3\} \{4\} \} , \{ \{1, 3\} \{2\} \{4\} \} , \\ & \{ \{1\} \{2, 3\} \{4\} \} , \{ \{1\} \{2\} \{3\} \{4\} \} \}\end{aligned}$$

$relabel(C_0)$  changes it to

$$\begin{aligned}relabel(C_0) = & \{ \{ \{1\} \{2, 3, 4\} \} , \{ \{1\} \{2, 3\} \{4\} \} , \{ \{1\} \{2, 4\} \{3\} \} , \\ & \{ \{1\} \{2\} \{3, 4\} \} , \{ \{1\} \{2\} \{3\} \{4\} \} \}\end{aligned}$$

The transformation to  $C_1$  yields  $C_1 = \{ \{ \{1\} \{2, 3\} \{4\} \} , \{ \{1\} \{2\} \{3\} \{4\} \} \}$ .  
The transformation to  $C_2$  yields  $C_2 = \{ \{ \{1\} \{2\} \{3\} \{4\} \} \}$ .

## 5 Practical implications

At the first glance Theorem 1 may seem to be trivial, useless, unrealistic and impractical. Trivial because it may appear to be obvious that a total-separation partition is closest to all the other partitions. Useless because nobody is interested in obtaining a consensus in terms of a total-separation partition. Unrealistic because the space of all possible clusterings is so immense, that for real world sample sizes one would never run so many clustering algorithms as to fill the whole partition universe. Impractical because one usually restricts the number of clusters  $k$  in a partition by an upper bound (much) lower than the number of elements  $n$ .

In order to demonstrate that these intuitions are wrong, we performed a series of simulation experiments results of which are summarized in Tables 1 (sampling the universe, no structure in the underlying data assumed), 2 (sampling the universe with a modified partition dissimilarity measure, no structure in the underlying data assumed), 3 (sampling a subuniverse where the presence of a simple structure in the data is assumed), and further ones, discussed in the next section.

With the experiments, we address the following questions:

- Does the bias of Rand index to choose the total-separation partition as consensus for the whole universe of partitions persist if we consider only a uniform random sample from this universe?
- Is this bias a general property of cluster quality indices or is it specific to Rand index?

- Does the bias of Rand index to choose the total-separation partition pertain if there is a structure in the partitions for which a consensus is sought?

While these questions are addressed in the current section, we pose them again in the next section in the context of constraining the set of partitions to those with an upper bound on the number of partitions.

Recall that consensus clustering uses Cluster Difference (derived from Rand Index) as a measure of distance between partitions in order to identify the consensus partition. We have already demonstrated theoretically, that if the set of our clustering algorithms would yield all possible partitions, then Rand index would pick up the total-separation partition. But of course the space of all partitions is too large so that we will never get all possible partitions. Nonetheless by manipulating clustering algorithm parameters in an irresponsible way we can obtain a random sample from this universe. In fact it is quite easy to invent clustering algorithms delivering for the same set of data any clustering we want. This section simulates such a situation and shows experimentally what the outcome of consensus clustering will be questionable also in such a case. See comments on experiments in tables 1 and 2 below. This suggests that the user exploiting the technology of consensus clustering must at least have an approximate vision of the geometry of the data space and parametrise clustering algorithms in a way not disturbing this geometry. Only in this case the consensus clustering may be helpful in the choice of appropriate compromise clustering. See the comments below on comparison of tables 1 and 3.

The experiments consisted in drawing 1,000 samples from the partition universe for each parameter setting (characterised by columns 1-4) and computed results are presented in column 5 (eventually column 6). So for example in Table 1 in the second data row we have the information, that samples were drawn from the universe of partitions over 4 elements, that the number of possible partitions is 15, out of them samples of 8 partitions were drawn which were intended to be 50% of the sample space, and the evaluation result was 88.2%.

The experiment underlying Table 1 was devised to find out how often the total-separation partition will turn out to be the consensus partition for a uniformly randomly drawn sample of partitions. The results are visible in column 5.

The experiment underlying Table 2 was essentially the same as the previous one, but instead of dissimilarity measure  $unCD$  its modification  $unCD_m(a, b) = unCD^{10}(a, b) = (unCD(a, b))^{10}$  was used. Again the total-separation partition occurrence as consensus partition was counted. The results are visible in column 5.

The experiment underlying Table 3 differs from the previous two in that not the whole universe of partitions of  $n$  elements was considered, but only those partitions for which there exists a "structure" that is where elements 1, 3 occur always in the same cluster. Here  $unCD$  was used as distance measure as in Table 1. The 5th column counts the number of occurrences of total-separation cluster as consensus, while the 6th column tells how frequently the expected partition (where all elements are singletons except for elements 1 and 3 that

constitute one cluster) is discovered as consensus partition.

Comparing Tables 1 and 2 (rows where the 3rd column contains 100% entry) we see first of all that the result of the Theorem 1 is *far from being trivial*. Total-separation partition does not need to be the default choice for a consensus of the entire universe of possible partitions. A skilled choice of dissimilarity function may point at any partition. It is the particular property of Rand Index (Cluster Difference) that distinguishes total-separation partition. So one can say that Rand Index is biased towards total-separation partition in case of no pattern in the data. This property seems to be vital and has never been reported in the context of consensus clustering. An investigation of other measures with respect to their behaviour under missing structure should be carried out.

Let us compare column 5 of tables 1 and 3. In table 1 we simulated the case that there was no intrinsic structure behind the data, so that the partitions obtained from various algorithms just provide random samples from the universe of all possible partitions. We count how frequently the total-separation partition will occur as the consensus between the diverse partitions. It turns out that it happens quite often even if we have small sets of partitions. In table 3 on the other hand such a situation never happens. Furthermore, in column 6 of this table one can see that the centre of the set of clusters exhibiting the assumed structure occurs quite often as the consensus in this experiment. So non-occurring of total-separation partition as a consensus under a sufficiently large set of competing partitions can be considered as a good indicator of existence of some structure in the data. So the ability of a consensus clustering algorithm to provide total-separation partition as the consensus is a very *useful* property because it allows us to discern between a meaningful and meaningless set of outcomes of diverse clustering algorithms.

The Theorem 1 provides us with an important insight into the partition space because the indicated behaviours are observed not only in the whole universe, but also in sufficiently big samples. Hence the distance of the actual consensus from the total-separation partition is an important indicator of the actual structure in the data. Enforcing exclusion of total-separation partition from consensus seeking algorithm is not a wise decision. So the result is very *practical*.

Let us also stress that for too small sample sizes one can get impression of existence of a structure in the data even if there is none. Hence in practice one should verify the validity of the consensus in the application domain.

We will postpone the discussion of the restriction of the number of clusters  $k$  in a partition by an upper bound to the next section as it requires some additional theoretical discussion.

In summary, the experiments allow to answer the posed questions as follows:

- The bias of Rand index to choose the total-separation partition as consensus for the whole universe of partitions *persists* if we consider only a uniform random sample from this universe?
- This bias *is not* a general property of cluster quality indices and the direction of the bias should be investigated separately for other cluster quality

Table 1: Results of consensus clustering of randomly selected partitions

data set size	number of all possible clusterings	sample size	percentage of possible clusterings	percentage of total-separation partitions in consensus
<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
4	15	15	100 %	100 %
4	15	8	50 %	88.2 %
5	52	52	100 %	100 %
5	52	26	50 %	100 %
5	52	21	40 %	97.2 %
5	52	10	20 %	78.9 %
5	52	5	10 %	23.1 %
6	203	203	100 %	100 %
6	203	102	50 %	100 %
6	203	81	40 %	100 %
6	203	41	20 %	99.9 %
6	203	20	10 %	96.5 %
7	877	88	10 %	100 %
7	877	44	5 %	99.9 %
7	877	35	4 %	99.7 %
7	877	18	2 %	93.7 %
7	877	9	1 %	49.9 %
8	4140	41	1 %	100 %
8	4140	21	0.5 %	94 %
8	4140	17	0.4 %	89 %
8	4140	8	0.2 %	61 %
8	4140	4	0.1 %	38 %

Table 2: Results of consensus clustering of randomly selected partitions for a modified distance measure

data set size	number of all possible clusterings	sample size	percentage of possible clusterings	percentage of total-separation partitions in consensus
<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
4	15	15	100 %	0 %
4	15	8	50 %	30.2 %
5	52	52	100 %	0 %
5	52	26	50 %	6.5 %
5	52	21	40 %	4.4 %
5	52	10	20 %	16.2 %
5	52	5	10 %	16.4 %
6	203	203	100 %	0 %
6	203	102	50 %	0.7 %
6	203	81	40 %	3.2 %
6	203	41	20 %	11.3 %
6	203	20	10 %	5.5 %
7	877	88	10 %	1.5 %
7	877	44	5 %	2.3 %
7	877	35	4 %	2.9 %
7	877	18	2 %	8.6 %
7	877	9	1 %	6.8 %

Table 3: Results of consensus clustering of randomly selected partitions from a set of partitions exhibiting simple structure

data set size	number of all cluster- ings with structure	sample size	percentage of possible clusterings	percentage of total- separation partitions in consensus	percentage of structure set centre
<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
4	5	5	100 %	0 %	100 %
4	5	2.5	50 %	0 %	69.2 %
5	15	15	100 %	0 %	100 %
5	15	8	50 %	0 %	89.7 %
5	15	6	40 %	0 %	74 %
5	15	3	20 %	0 %	22.8 %
5	15	2	10 %	0 %	53.1 %
6	52	52	100 %	0 %	100 %
6	52	26	50 %	0 %	100 %
6	52	21	40 %	0 %	97.2 %
6	52	10	20 %	0 %	77.3 %
6	52	5	10 %	0 %	23.5 %
7	203	41	20 %	0 %	99.9 %
7	203	20	10 %	0 %	95.8 %
7	203	16	8 %	0 %	90.4 %
7	203	8	4 %	0 %	70.2 %
7	203	4	2 %	0 %	47.6 %



indices.

- Rand index behaves differently if there is a structure in the partitions for which a consensus is sought.

## 6 Imposing a limit on the number of clusters to be considered

Partitions explored practically in the case of consensus clustering always contain much fewer clusters compared to the size of the set of objects ( $k \ll n$ ). The question is: would a limitation on the number of clusters  $k < k_{max}$  change anything with respect to the previous results?

Note that for a partition consisting of  $k$  sets each of its complex extensions contains also  $k$  sets, but the simple extension contains  $k + 1$  sets. Therefore, as long as  $k < k_{max}$ , nothing changes in the discussion of sections 4.4 and 4.5. But for  $k = k_{max}$ , we need to update equations (6), (7) and (8), because we will no longer consider (hence count) distances to the simple extensions. So the equation (6) has to be substituted by

$$\begin{aligned}
\sum_{l=1, l \neq l'}^k unCD(\Gamma_l, \Gamma_{l'}) &= \sum_{l=1, l \neq l'}^k (card(S_l - S_{l'}) + card(S_{l'} - S_l)) \\
&= \sum_{l=1, l \neq l'}^k (card(S_l) + card(S_{l'})) \\
&= (k - 1) \cdot card(S_{l'}) + (n - card(S_{l'})) \\
&= (k - 2) \cdot card(S_{l'}) + n \geq n
\end{aligned} \tag{9}$$

which holds of course only if  $k_{max} > 1$ , which is rather a non-restrictive assumption.

The equation (7) has to be substituted by

$$\begin{aligned}
\sum_{l=1}^k unCD(\Gamma_0, \Gamma'_l) &= \sum_{l=1}^k (unCD(\Gamma^*, \Gamma^{*'}; n) + card(S'_l)) \\
&= n + \sum_{l=1}^k unCD(\Gamma^*, \Gamma^{*'}; n) \\
&= n + k \cdot unCD(\Gamma^*, \Gamma^{*'}; n)
\end{aligned} \tag{10}$$

The equation (8) has to be substituted by

$$\begin{aligned}
\sum_{l=1}^k unCD(\Gamma_m, \Gamma'_l) &= \sum_{l=1}^k (unCD(\Gamma^*, \Gamma^{*'}; n) + card(S_m - S'_l) + card(S'_l - S_m)) \\
&= \sum_{l=1}^k (unCD(\Gamma^*, \Gamma^{*'}; n) + card(S_m) - card(S_m \cap S'_l) \\
&\quad + card(S'_l) - card(S'_l \cap S_m)) \\
&= \sum_{l=1}^k unCD(\Gamma^*, \Gamma^{*'}; n) + k \cdot card(S_m) + n \\
&\quad - \sum_{l=1}^k card(S_m \cap S'_l) - \sum_{l=1}^k card(S'_l \cap S_m) \\
&= k \cdot unCD(\Gamma^*, \Gamma^{*'}; n) + k \cdot card(S_m) + n \\
&\quad - card(S_m) - card(S_m) \\
&= k \cdot unCD(\Gamma^*, \Gamma^{*'}; n) + (k - 2) \cdot card(S_m) + n \\
&\geq k \cdot unCD(\Gamma^*, \Gamma^{*'}; n) + n
\end{aligned} \tag{11}$$

which is again valid under  $k_{max} > 1$ .

Hence lemmas 1 and 2 retain their validity for the restricting  $k_{max}$  and hence the reasoning presented in section 4.6.

So the following theorem holds:

**Theorem 2.** *For any  $n$  objects, among all partitions the partition where each object falls into a separate cluster (called subsequently total-separation partition) has the lowest average Cluster Difference  $CD$  to the other partitions consisting of at most  $k_{max} > 1$  clusters.*

But of course the limitation of the number of clusters still leaves a huge space of possible partitions of the set of elements. So let us return to the discussion of sampling this space.

Experiments analogous to those of previous section have been performed and are summarized in Table 4. Random samples from the space of permissible clusterings (with  $k < k_{max}$ ) were drawn and the suitability of total-separation partition as consensus cluster was checked. A similar pattern to table 1 was observed – with sufficiently large samples total-separation partition indicates that there is no real relationship underlying the various clusterings.

The question can be raised: Can it really be so bad that modern day clustering techniques would provide a clustering when no clusters are there in the data. Vast majority of clustering algorithms produce partitions whatever data they get. We can just point at  $k$ -means algorithm, but others, like DBSCAN, single-link etc. could be used. Imagine a large collection of data points in a high-dimensional space. Furthermore imagine the points are randomly uniformly distributed in space. Imagine that in order to perform  $k$ -means clustering

Table 4: Results of consensus clustering of randomly selected partitions with upper limit on  $k$  equal 4

data set size	number of all possible clusterings	sample size	percentage of possible clusterings	percentage of total-separation partitions in consensus
<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
6	187	187	100 %	100 %
6	187	94	50 %	100 %
6	187	75	40 %	100 %
6	187	37	20 %	99.6 %
6	187	19	10 %	81.5 %
7	715	143	20 %	100 %
7	715	72	10 %	100 %
7	715	57	8 %	99.8 %
7	715	29	4 %	95.5 %
7	715	14	2 %	76.6 %

efficiently, samples from this collection are drawn and  $k$ -means algorithm is performed on each of them. As  $k$ -means always clusters the entire space (computes centroids), we would have in fact a random sample of partitions from this universe. Now let the many partitions be processed by consensus clustering using Rand Index. A consensus partition will be undoubtedly found. The lesson from our theorem will be that it should converge towards total-separation partition.

But we can face also another serious real problem. Assume that we have a data set with an underlying structure (say with two dense areas separated by empty space), but located in a highly dimensional space.  $k$ -means algorithm has the known property of being  $k$ -rich that is upon proper transformation of distances between data points ANY clustering consisting of  $k$  clusters may be obtained. And upon applying various clustering algorithms indeed distances are transformed, e.g. by standardization, normalization, spectral transformations etc. This means, however, that upon irresponsible choice of distance transformations we may sample from the universe of all possible partitions in spite of the fact that the data has originally a structure. And here again, as a result of consensus clustering, we can get total-separation partition, which would be really bothering. So any transformation we apply to the data as an element of the clustering process must not violate the geometrical structure of the data space we expect to see.

So whether or not we restrict ourselves to an upper limit of clusters in a partition, same answers are to be given to our questions driving the experiments. So either we have to live with the risk that a consensus clustering returns us a meaningless consensus or have to further develop these consensus methods so that they are able to refuse to return a partition if no real structure in the data

exists.

## 7 Remarks on meta-clustering

The formulas derived in the preceding sections shed also some light on possible outcome of meta-clustering. It is obvious that assuming a sufficiently "large" number of clusters, the simple extensions will be cluster centres and clusters will consist of extensions of the same reduct.

Let us introduce the concept of  $p$ -th order reduct and  $p$ -th order extension. If  $\Gamma^*$  is a reduct of  $\Gamma$ , then it is the 1st order reduct of  $\Gamma$ , and  $\Gamma$  is the 1st order extension of  $\Gamma^*$ . If it is a simple extension, then it is the 1st order simple extension. If  $\Gamma^+$  is a  $p$ -th order reduct of  $\Gamma^*$  and  $\Gamma^*$  is a reduct of  $\Gamma$ , then  $\Gamma^+$  is the  $(p+1)$ st order reduct of  $\Gamma$  and  $\Gamma$  is the  $(p+1)$ st order extension of  $\Gamma^+$ . If at the same time  $\Gamma^*$  is a  $p$ -th order simple extension of  $\Gamma^+$  and  $\Gamma$  is a simple extension of  $\Gamma^*$ , then  $\Gamma$  is the  $(p+1)$ st order simple reduct of  $\Gamma^+$ . Otherwise  $\Gamma$  is its complex extension.

With the above definition let us arrange a hierarchical clustering, where the  $p$ -th level (from the bottom) consists of clusters with common  $p$ -th degree reduct. In this case it is easy to see from Theorem 1 that the cluster centres would be  $p$ -th degree simple extensions of the respective reduct, around which the cluster is defined. This hierarchy would be a local minimum for the combined distance of elements from their cluster centres at respective levels (each element of a cluster is closer to its own cluster centre than to any other cluster centre). One can in fact check that at the top level of such a hierarchy moving any object between classes is not possible.

Let us demonstrate it by considering two meta-cluster centres over the set of partitions of  $n$  elements:  $\Gamma_0$  being total-separation partition and  $\Gamma_1$  in which elements 1, 2 are in one cluster and the others in separate clusters. Consider now a partition  $\Gamma_q$  in which elements 1 and 2 are in separate clusters, and no assumption with respect to other is done.

$$\begin{aligned} unCD(\Gamma_0, \Gamma_q) &= I_{12/0q} + \sum_{j=3}^n I_{1j/0q} + \sum_{j=3}^n I_{2j/0q} + \sum_{i=3}^{n-1} \sum_{j=i+1}^n I_{ij/0q} \\ unCD(\Gamma_1, \Gamma_q) &= I_{12/1q} + \sum_{j=3}^n I_{1j/1q} + \sum_{j=3}^n I_{2j/1q} + \sum_{i=3}^{n-1} \sum_{j=i+1}^n I_{ij/1q} \end{aligned} \quad (12)$$

where  $I_{jk/lm}$  means indicator of membership of element  $j, k$  in same cluster in one of partitions  $\Gamma_l, \Gamma_m$  and in different in the other. It is easily checked that  $\sum_{j=3}^n I_{1j/0p} = \sum_{j=3}^n I_{1j/1p}$ ,  $\sum_{j=3}^n I_{2j/1p} = \sum_{j=3}^n I_{2j/0p}$  and  $\sum_{i=3}^{n-1} \sum_{j=i+1}^n I_{ij/0p} = \sum_{i=3}^{n-1} \sum_{j=i+1}^n I_{ij/1p}$ . The only difference is  $I_{12/1q}$  equal to one and  $I_{12/0q}$  equal to zero.  $\Gamma_q$  is closer to  $\Gamma_0$  than to  $\Gamma_1$ , as expected. If on the other hand both 1 and 2 would be in the same cluster in  $\Gamma_q$ , the situation would be inverted.

Let us consider more detailed levels of meta-clustering, that is let  $\Gamma_0$  and  $\Gamma_1$  be  $p$ th simple extensions of  $\Gamma_0^*$  and  $\Gamma_1^*$  (over same set of elements) resp. and let  $\Gamma_q$  be  $p$ th extension of  $\Gamma_0^*$ .

$$\begin{aligned} unCD(\Gamma_0, \Gamma_q) &= \sum_{i=1}^{n-p} \sum_{j=i+1}^{n-p} I_{ij/0q} + \sum_{i=1}^{n-p} \sum_{j=n-p+1}^n I_{ij/0q} + \sum_{i=n-p+1}^{n-1} \sum_{j=i+1}^n I_{ij/0q} \\ unCD(\Gamma_1, \Gamma_q) &= \sum_{i=1}^{n-p} \sum_{j=i+1}^{n-p} I_{ij/1q} + \sum_{i=1}^{n-p} \sum_{j=n-p+1}^n I_{ij/1q} + \sum_{i=n-p+1}^{n-1} \sum_{j=i+1}^n I_{ij/1q} \end{aligned} \quad (13)$$

Again obviously  $\sum_{i=1}^{n-p} \sum_{j=n-p+1}^n I_{ij/0q} = \sum_{i=1}^{n-p} \sum_{j=n-p+1}^n I_{ij/1q}$ ,  $\sum_{i=n-p+1}^{n-1} \sum_{j=i+1}^n I_{ij/0q} = \sum_{i=n-p+1}^{n-1} \sum_{j=i+1}^n I_{ij/1q}$ , so that again  $\sum_{i=1}^{n-p} \sum_{j=i+1}^{n-p} I_{ij/0q} = 0$  (as both  $\Gamma_0$  and  $\Gamma_q$  have the same reduct) and  $\sum_{i=1}^{n-p} \sum_{j=i+1}^{n-p} I_{ij/1q} > 0$  make the difference. One concludes that if the meta-clustering is a hierarchical one and the clusters are build around same reducts then the clusters at each hierarchy level are stable.

Hence it is obvious that also meta-clustering is biased - there exists a structure in spite of filling in the whole sample space.

The above meta-clustering "algorithm" was pretty much manual. One can ask whether or not other algorithms will exhibit the same tendency. In particular if it turns out that the total-separation partition is chosen as centre of any of the meta-clusters.

Let us try out  $k$ -medoids clustering, implemented in  $R$  as *pam* algorithm and the popular  $k$ -means algorithm as meta-clustering methods. For this purpose let us span a  $(n-1) \cdot n/2$  dimensional space for partitions of  $n$  elements. For a partition  $\Gamma$ , for  $1 \leq i < j \leq n$  the  $(j-1) \cdot (j-2)/2 + i$ -th coordinate in this space would be equal to 1 if both elements  $i, j$  belong to the same cluster and equal to 0 otherwise. It is easily seen that  $unCD$  is the taxicab-distance between partitions in this space or the square of Euclidean distance, making it reasonable to apply e.g.  $k$ -means algorithm. Let us have a look at the meta-clusters generated using  $k$ -means and *pam* algorithm (as implemented in  $R$  system) under this representation.

For a set of 4 elements,  $k$ -means algorithm with  $k = 2$  splits the set of all partitions into the meta-clusters: Meta-cluster 1 with minimal distance to meta-cluster centre 1.5 containing 10 partitions:

$\{\{1, 2, 4\}\{3\}\},$   
 $\{1, 2\}\{3, 4\}\},$   
 $\{\{1, 2\}\{3\}\{4\}\},$   
 $\{\{1, 4\}\{2, 3\}\},$   
 $\{\{1\}\{2, 3, 4\}\},$   
 $\{\{1\}\{2, 3\}\{4\}\},$   
 $\{\{1, 4\}\{2\}\{3\}\},$   
 $\{\{1\}\{2, 4\}\{3\}\},$

$\{\{1\}\{2\}\{3,4\}\},$   
 $\{\{1\}\{2\}\{3\}\{4\}\}.$

The last one is the closest point to centroid of this meta-cluster. Meta-cluster 2 with minimal distance to meta-cluster centre 2 containing 5 partitions:

$\{\{1,2,3,4\}\},$   
 $\{\{1,2,3\}\{4\}\},$   
 $\{\{1,3,4\}\{2\}\},$   
 $\{\{1,3\}\{2,4\}\},$   
 $\{\{1,3\}\{2\}\{4\}\}$

The last one is the closest point to centroid of this meta-cluster.

On the other hand, the *pam* algorithm with split parameter set to 2 creates a meta-cluster consisting of partitions  $\{\{1,2,3,4\}\}$  and  $\{\{1\}\{2,3,4\}\}$ , while the other meta-cluster contains all the other partitions, with  $\{\{1\}\{2\}\{3\}\{4\}\}$  (total-separation partition) being the medoid.

For 5 elements *k-means* provides two meta-clusters: meta-cluster 1 minimum distance to meta-cluster centroid equal 2.432 and of card. 37 with element closest to the centroid  $\{\{1\}\{2\}\{3\}\{4\}\{5\}\}$  and a meta-cluster 2 (min. dist. 3 card. 15) with element closest to the centroid  $\{\{1,2\}\{3\}\{4\}\{5\}\}$ .

The two meta-clusters returned by *pam* have cardinality 6 with medoid  $\{\{1,2,3,4,5\}\}$  and card. 46 with medoid  $\{\{1,2\}\{3\}\{4\}\{5\}\}$ . Here, total-separation partition was not selected, but still is among the candidates.

For 6 elements *k-means* creates one meta-cluster of card. 156 with the total-separation partition being the closest one to the centroid, while the other meta-cluster of card. 47 has 25 elements with minimal distance of 7.40 from the centroid.

*pam* provides two meta-clusters of card. 52 and 151 with medoids  $\{\{1\}\{2\}\{3\}\{4\}\{5,6\}\}$  and total-separation one resp.

For 7 elements *k-means* returns a meta-cluster of card. 582 with the total-separation partition being the closest one (dist. 3.463) to the centroid, and another meta-cluster of card. 295 105 elements closest to the centroid (dist. 9.264406)

*pam* provides two meta-clusters of card. 203 and 674 with medoids  $\{\{1\}\{2\}\{3\}\{4\}\{5\}\{6,7\}\}$  and total-separation one resp.

For 8 elements in *k-means* we have one meta-cluster of card. 2570 with the total-separation partition being the closest one to the centroid (dist. 4.21), while the other meta-cluster of card. 1570 has 700 elements with minimal distance of 11.37 from the centroid.

*pam* provides two meta-clusters of card. 877 and 3263 with medoids  $\{\{1\}\{2\}\{3\}\{4\}\{5\}\{6\}\{7,8\}\}$  and total-separation one resp.

For *k-means* it should be underlined that the distance between cluster centres does not exceed 1 in any of the above cases.

It is worth noting that the meta-cluster around the total-separation partition is the larger one. In fact, with *k-means* for 8 elements if *k* grows even to 100, this meta-cluster (with total-separation partition) is the largest one. The percentage of variance explained (betweenness/ totals) is however low. With

k=2 3% is explained, with k=20 26% is explained, with k=50 41% is explained, with k=100 54%, with k=150 60%, with k=200 64%.

We also see in these experiments that split into two classes as performed by *pam* is in agreement with our claim that  $p$ -th order simple extensions will become meta-cluster centres. This is even true for *k-means* with the total separation partition.

A comment on the difference between *pam* and *k-means* results seems to be required. First of all we made here a trick of pressing the distances between partitions into a vector space. This should not affect *pam* algorithm as the distances in this space are the same as the original space, but it may be a little bit confusing for *k-means* which actually uses the square roots of the original distances. Furthermore, the dense partition space was replaced by a sparse vector space. This leads *k-means* to explore places in this space (as centroids) that also not have a clear interpretation. Possibly one could investigate a kind of fuzzification here in order to have an insight into what the centroid may mean. This may be a future research path.

Finally, let us have a look at the performance of *pam* when we do not (meta)-cluster the universe of all possible partitions, but only random samples from it. We check again if the total-separation partition is among the candidates. We have, however, one difference here compared to the settings of experiments in Section 5. As *pam* seeks for medoids, total-separation partition must be among the partitions to be meta-clustered, so we enforce selecting it when randomly sampling. Table 5 shows the behaviour of *pam* when the samples are uniformly drawn from the universe of partitions. Total-separation partition is the dominating option for medoid of at least one of two meta-clusters, that *pam* is requested to generate. Table 6 shows the behaviour of *pam* when we draw the samples from the sub-universe with a structure as done in Table 3. Here, for sufficiently large random samples the total-separation partition, though present in each sample, is rarely chosen.

## 8 Conclusions and future work

As the number of available clustering algorithms applicable to the same data is growing, and the potential outputs may differ substantially, methodologies to reconcile them like meta-clustering or consensus clustering are under development.

In this paper we demonstrated that both consensus clustering and meta-clustering using Cluster Difference (derived from Rand Index) as a measure of distance between partitions, when applied to the universe of all possible partitions, point to the partition containing each element in a separate set as the best compromise.

This suggests that the user performing the task of clustering must at least have an approximate vision of the geometry of the data space. Only in this case the mentioned techniques may be helpful in the choice of appropriate compromise clustering.

Table 5: Results of meta-clustering via *pam* of randomly selected partitions

data set size	number of all possible clusterings	sample size	percentage of possible clusterings	percentage of total-separation partitions as meta-cluster centre
<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
4	15	15	100 %	100 %
4	15	8	50 %	64.9 %
5	52	52	100 %	100 %
5	52	27	50 %	99.6 %
5	52	21	40 %	97.8 %
5	52	11	20 %	83.7 %
5	52	6	10 %	68.9 %
6	203	203	100 %	100 %
6	203	102	50 %	100 %
6	203	82	40 %	100 %
6	203	41	20 %	99.9 %
6	203	21	10 %	98 %
7	877	89	10 %	100 %
7	877	45	5 %	100 %
7	877	36	4 %	100 %
7	877	19	2 %	99.6 %
7	877	10	1 %	96.9 %
8	4140	84	2 %	100 %
8	4140	42	1 %	100 %
8	4140	34	0.8 %	100 %
8	4140	18	0.4 %	99 %
8	4140	9	0.2 %	99 %



Table 6: Results of meta-clustering via *pam* of randomly selected partitions out of a set with structure

data set size	number of all possible clusterings	sample size	percentage of possible clusterings	percentage of total-separation partitions as meta-cluster centre
<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
4	6	6	100 %	0 %
4	6	3	50 %	10.8 %
5	16	16	100 %	0 %
5	16	9	50 %	0.2 %
5	16	7	40 %	1.4 %
5	16	4	20 %	10.1 %
5	16	3	10 %	25 %
6	53	53	100 %	0 %
6	53	27	50 %	0 %
6	53	22	40 %	0 %
6	53	11	20 %	1.5 %
6	53	6	10 %	9.5 %
7	204	42	20 %	0.1 %
7	204	21	10 %	3.6 %
7	204	17	8 %	5.1 %
7	204	9	4 %	19.6 %
7	204	5	2 %	35.2 %
8	878	71	8 %	6 %
8	878	36	4 %	16 %
8	878	29	3.2 %	22 %
8	878	15	1.6 %	32 %
8	878	8	0.8 %	43 %

It seems also worth investigating, how other cluster quality functions used as distances between partitions would behave under consensus clustering of the space of all possible partitions.

It seems also worth investigating how such measures would behave not in the full universe of all partitions but rather for uniform random samples of it. Such a sampling would then constitute a background for investigations into the behaviour of other partition comparison indexes, of consensus and meta-clustering methods as well as for checking if a resultant consensus-partition or meta-cluster really gives a new insight or is just a random artefact.

## Acknowledgements

The author wishes to thank to the Institute of Computer Science of Polish Academy of Sciences for promoting and financing this research.

## References

- [1] Anderberg, M.: *Cluster Analysis for Applications*, Academic Press, London, 1973.
- [2] Barthelemy, J.-P., Leclerc, B.: The median procedure for partition, *Partitioning Data Sets* (I. C. et al, Ed.), AMS DIMACS Series in Discrete Mathematics, 1995.
- [3] Bifulco, I., Fedullo, C., Napolitano, F., Raiconi, G., Tagliaferri, R.: Multiple data structure discovery through global optimisation, meta clustering and consensus methods, *International Journal of Knowledge Engineering and Soft Data Paradigms, v.1 n.4, October 2009*, 2009.
- [4] Bifulco, I., Iorio, F., Napolitano, F., Raiconi, G., Tagliaferri, R.: Interactive Visualization Tools for Meta-Clustering, *Proceedings of the 2009 conference on New Directions in Neural Networks: 18th Italian Workshop on Neural Networks: WIRN 2008*, IOS Press, Amsterdam, The Netherlands, The Netherlands, 2009, ISBN 978-1-58603-984-4.
- [5] Caruana, R., Elhawary, M., Nguyen, N., Smith, C.: Meta Clustering, *Proceedings of the Sixth International Conference on Data Mining, ICDM '06*, IEEE Computer Society, Washington, DC, USA, 2006, ISBN 0-7695-2701-9.
- [6] Cui, Y., Fern, X. Z., Dy, J. G.: Learning multiple nonredundant clusterings, *ACM Transactions on Knowledge Discovery from Data (TKDD)*, **4**, October 2010, 15:1–15:32, ISSN 1556-4681.
- [7] Dasgupta, S., Ng, V.: Which clustering do you want? inducing your ideal clustering with minimal feedback, *J. Artif. Int. Res.*, **39**, September 2010, 581–632, ISSN 1076-9757.

- [8] Fowlkes, E. B., Mallows, C. L.: .
- [9] Ghosh, J., Acharya, A.: Cluster ensembles, *Wiley Interdisc. Rev.: Data Mining and Knowledge Discovery*, **1**(4), 2011, 305–315.
- [10] Gionis, A., Mannila, H., Tsaparas, P.: Clustering Aggregation, *ACM Trans. Knowl. Discov. Data*, **1**(1), March 2007, ISSN 1556-4681.
- [11] Goder, A., Filkov, V.: Consensus Clustering Algorithms: Comparison and Refinement, *Proceedings of the Workshop on Algorithm Engineering and Experiments, ALENEX 2008, San Francisco, California, USA, January 19, 2008* (J. I. Munro, D. Wagner, Eds.), 2008.
- [12] Goder, A., Filkov, V.: Consensus Clustering Algorithms: Comparison and Refinement., *Alenex*, 8, SIAM, 2008.
- [13] Gordon, A., Vichi, M.: Partitions of partitions, *Journal of Classification*, **15**, 1998, 265–285.
- [14] Hore, P., Hall, L. O., Goldgof, D. B.: A scalable framework for cluster ensembles, *Pattern Recogn.*, **42**(5), May 2009, 676–688, ISSN 0031-3203.
- [15] Lei, Y., Bezdek, J. C., Romano, S., Vinh, N. X., Chan, J., Bailey, J.: Ground Truth Bias in External Cluster Validity Indices, *CoRR*, **abs/1606.05596**, 2016.
- [16] Li, T., Ding, C.: Weighted consensus clustering, *Proceedings of 2008 SIAM International Conference on Data Mining (SDM 2008), Atlanta, April 24-26, 2008*, Society for Industrial and Applied Mathematics, 2008.
- [17] Milligan, G. W., Cooper, M. C.: A study of the comparability of external criteria for hierarchical cluster analysis, *Multivar. Behav. Res.*
- [18] Monti, S., Tamayo, P., Mesirov, J., Golub, T.: Consensus Clustering: A Resampling-Based Method for Class Discovery and Visualization of Gene Expression Microarray Data, *Mach. Learn.*, **52**(1-2), July 2003, 91–118, ISSN 0885-6125.
- [19] Morlini, I., Zani, S.: Comparing Approaches for Clustering Mixed Mode Data: An Application in Marketing Research, *Data Analysis and Classification: Proceedings of the 6th Conference* (F. Palumbo, C. N. Lauro, M. Greenacre, Eds.), Springer, 2010.
- [20] Nguyen, N., Caruana, R.: Consensus Clusterings, *Proceedings of the 7th IEEE International Conference on Data Mining (ICDM 2007), October 28-31, 2007, Omaha, Nebraska, USA*, 2007.
- [21] Niu, D., Dy, J. G., Jordan, M. I.: Multiple Non-Redundant Spectral Clustering Views., *ICML ’10*, 2010.

- [22] Punera, K., Ghosh, J.: Consensus Based Ensembles of Soft Clusterings, *Applied Artificial Intelligence: An International Journal*, **22**(7-8), 2008, 780–810.
- [23] Simone Romano, James Bailey, V. N. K. V.: Standardized Mutual Information for Clustering Comparisons: One Step Further in Adjustment for Chance, *Proceedings of The 31st International Conference on Machine Learning*, 2014.
- [24] Strehl, A., Ghosh, J.: Cluster ensembles — a knowledge reuse framework for combining multiple partitions, *J. Mach. Learn. Res.*, **3**, March 2003, 583–617, ISSN 1532-4435.
- [25] Topchy, A., Jain, A. K., Punch, W.: Clustering ensembles: Models of consensus and weak partitions, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **27**, 2005, 1866–1881.
- [26] Vogel, T., Naumann, F.: Semi-Supervised Consensus Clustering: Reducing Human Effort, *Proceedings of the International Workshop on Data Integration and Applications*, 2014.
- [27] Wang, H., Shan, H., Banerjee, A.: Bayesian cluster ensembles, *Statistical Analysis and Data Mining*, **4**, 2011, 54–70.
- [28] Wang, Y., Pan, Y.: Semi-Supervised Consensus Clustering for Gene Expression Data Analysis, *BioData Mining*, **7**(7), 2014.